
Descoperirea cunoștințelor conceptuale în baze de date

REZUMAT AL TEZEI DE DOCTORAT



Katalin Tünde Jánosi-Rancz

Școala Doctorală de Matematică și Informatică

Conducător:

Prof. Dr. Zoltán Kása

Departmentul de Informatică
Facultatea de Matematică și Informatică
Universitatea Babeș-Bolyai
Cluj-Napoca, România

2014

Lista publicațiilor

1. **K.T. Janosi-Rancz**, A. Lajos, Semantic Data Extraction, Accepted for The 8th International Conference INTER-ENG 2014, Interdisciplinarity in Engineering, 9 - 10 October 2014, "Petru Maior" University of Tîrgu Mureș, Romania
2. **K.T. Janosi-Rancz**, Finding, Managing and Inforcing CFDs and Ars via a semi-automatic learning strategy, 10th Joint Conference on Mathematics and Computer Science, May 21-25, Cluj Napoca, Romania, 2014
3. **K.T. Janosi-Rancz**, V. Varga, XML Schema Refinement Through Formal Concept Analysis, Studia Universitatis Babes-Bolyai, Informatica, Volume LVII, Number 3, 2012, 49-64
4. V. Varga, **K. T. Janosi Rancz**, C. Sacarea, K. Csioban, XML Design: an FCA Point of View, Proceedings of 2010 IEEE International Conference on Automation, Quality and Testing, Robotics (AQTR), Theta 17th edition, 2010, 165-170
5. **K. T. Janosi-Rancz**, V. Varga and T. Nagy, Detecting XML Functional Dependencies through Formal Concept Analysis, 14th East European Conference on Advances in Databases and Information Systems (ADBIS), Novi Sad, Serbia, Springer LNCS 6295, 2010, 595-598 (**Citat de 2 ori**)
6. D. Ignatov, **K. T. Janosi-Rancz**, S. Kuznetsov, Towards a framework for near-duplicate detection in document collections based on closed sets of attributes. Acta Universitatis Sapientiae, Informatica, 2/2009, 215-233
7. **K. T. Janosi Rancz**, V. Varga, J. Puskas, A Soft Tool for Relational Database Design using Formal Concept Analysis, 7th Joint Conference on Mathematics and Computer Science, Studia Universitatis Babes-Bolyai, Informatica, vol. LIII, No. 2, 2008, 67-78 (**Citat de 1**)
8. **K. T. Janosi Rancz**, V. Varga, A method for mining functional dependencies in

relational database design using FCA. *Studia Universitatis Babes-Bolyai, Informatica*, vol. LIII, No. 1, 2008, 17-28 (**Citat de 6 ori**)

9. V.Varga, **K. T. Janosi Rancz**, A Software Tool to Transform Relational Databases in Order to Mine Functional Dependencies in it Using Formal Concept Analysis, Proceedings of the Sixth Int. Conf. on Concept Lattices and Their Applications (CLA), Olomouc, Czech Republic, 2008, 1-9 (**Citat de 2 ori**)

Cuvinte cheie: Knowledge Discovery, XML design, Formal Concept Analysis, Semantic Search, Semantic Set, Data Mining, Information Extraction, Dependency Mining, Association Rule, Dependency Management, Concept Lattice, Data Management, Optimization, Data Extraction, Information retrieval

Cuprinsul tezei

1	Introducere	1
1.1	Descoperirea cunoștințelor conceptuale în baze de date	1
1.2	Contribuții și structura tezei	3
1.3	Contribuții publicate anterior	5
2	Fundamentele teoretice	7
2.1	Formal Concept Analysis	7
2.1.1	Conceptele de bază ale FCA	7
2.1.2	Context, concept și laticea conceptuală	9
2.1.3	Scalarea conceptuală	12
2.1.4	Literatura FCA	13
2.2	Dependențele în baze de date	13
2.2.1	Dependențele în baze de date relaționale	13
2.2.2	Dependențele în XML	15
3	Caracterizarea dependențelor funcționale (FD) cu contexte formale	25
3.1	Introducere	25
3.2	FD-uri calculate cu FCA	26
3.3	Crearea contextului formal folosind fișiere Inverted Index	29
3.3.1	Experimente	33
3.4	FCAFuncDepMine	36
3.5	Concluzii	38
4	Dependențe funcționale XML cu FCA	41
4.1	Introducere	41
4.2	Mineritul dependențelor funcționale XML folosind FCA	42
4.2.1	Prezentarea abordării	44
4.2.2	Crearea laticii conceptuale	46
4.2.3	Procesarea rezultatului FCA	47
4.2.4	Mineritul XFD-urilor conform ierarhia conceptuală	47
4.2.5	Găsirea cheilor XML	50

4.2.6	Detectarea redundanțelor în XML	50
4.2.7	Normalizare	51
4.3	Concluzii	51
5	Managementul dependențelor și prevenirea incoerențelor	53
5.1	Introducere	54
5.2	Învățarea, prezentarea, validarea și uitarea	56
5.2.1	Învățarea	56
5.2.2	Utilitatea	59
5.2.3	Validarea	60
5.2.4	Uitarea	62
5.3	Rezultate experimentale	64
5.4	Studierea SD-urilor cu FCA	65
5.4.1	Analizarea contextului formal pentru descoperirea cunoștințelor	68
5.5	Concluzii și cercetări viitoare	70
6	Extragerea datelor prin utilizarea arborilor semantici	73
6.1	Introducere	73
6.2	Descrierea metodei noastre	77
6.2.1	Limbajul <i>RK</i>	77
6.2.2	Arborele semantic	81
6.2.3	Extragerea cu <i>RK</i>	84
6.2.4	Navigarea cu <i>RK</i>	86
6.2.5	<i>RK</i> și schimbarea structurală	86
6.3	Rezultate experimentale	87
6.3.1	Compararea cu metode existente	89
6.4	Concluzii și cercetări viitoare	91
7	Concluzii și cercetări viitoare	93
7.1	Sumar	93
7.2	Perspectivă	96
	Bibliografie	99

Conținutul rezumatului

1	Introducere	1
1.1	Descoperirea cunoștințelor conceptuale în baze de date	1
1.2	Contribuții și structura tezei	2
1.3	Contribuții publicate anterior	4
2	Fundamentele teoretice	5
3	Caracterizarea dependențelor funcționale (FD) cu contexte formale	7
3.1	Rezultatele experimentului	8
4	Dependențe funcționale XML cu FCA	11
5	Managementul dependențelor și prevenirea incoerențelor	15
5.1	Rezultatele experimentului	16
6	Extragerea datelor prin utilizarea arborilor semantici	19
6.1	Rezultatele experimentului	21
7	Concluzii și planuri de viitor	23
	Bibliografie	25

1

Introducere

1.1 Descoperirea cunoștințelor conceptuale în baze de date

Odată cu creșterea de fără precedent a generării și acumulării datelor on-line, a crescut și dependența noastră de baze de date în privința stocării acestor date. Există o necesitate urgentă pentru noi metode de abordare a acestei provocări de lungă durată. Cum putem identifica, extrage și prezenta în mod eficient informațiile utile care sunt incluse în aceste baze de date? Și mai ales cum putem să le prezentăm într-un mod benefic nevoilor actuale și viitoare ale utilizatorului final? Găsirea răspunsurilor la aceste întrebări a fost sursa noastră de inspirație pentru a dezvolta abordările care sunt prezentate în această teză.

Descoperirea cunoștințelor în baze de date (Knowledge Discovery in Databases - KDD) [35] are ca scop dezvoltarea unor metode, tehnici și instrumente care ajută analiștii umani în procesul general de a obține unități de cunoaștere care la rândul lor pot fi utilizate în rezolvarea problemelor reale. Conceptele sunt necesare pentru a exprima cunoașterea umană [81]. În mod recent teoria laticelor a dus în gândirea matematică la reprezentarea cunoașterii și descoperirea cunoștințelor. Analiza formală a conceptelor (Formal Concept Analysis - FCA) [80; 26; 7; 10; 58; 44] oferă o abordare naturală pentru a descoperi concepte formale în datele care sunt descrise în forma unui context formal, descoperirea dependențelor de date și vizualizarea acestora printr-o singură structură conceptuală numită laticea conceptuală. Reprezentarea grafică a laticelor conceptuale s-a dovedit a fi extrem de utilă în descoperirea și înțelegerea relațiilor conceptuale într-un anumit set de date. Laticea conceptuală s-a dovedit a fi utilă și în alte aplicații referitoare la prelucrarea informațiilor și a cunoștințelor cum ar fi vizualizarea, analiza datelor, mineritul datelor (data mining) și managementul cunoștințelor.

FCA asigură suport în așa numita descoperire a cunoștințelor conceptuale în baze de date (Conceptual Knowledge Discovery in Databases - CKDD) care are ca scop facilitarea unui proces centrat pe om, de a descoperi cunoștințe în date prin vizualizarea și analiza structurii conceptuale a datelor [39].

Scopul nostru a fost de a extrage cunoștințe conceptuale de înalt nivel din baze de date. Contribuția noastră principală este mineritul diferitelor tipuri de dependențe de date (în date relaționale și XML) utilizând FCA; ne concentrăm, de asemenea, pe managementul dependențelor, prevenirea incoerențelor și extragerea de date semantice. Toate cercetările noastre sunt efectuate în domeniul CKDD.

1.2 Contribuții și structura tezei

În această secțiune vom prezenta principalele noastre contribuții și modul în care acestea sunt incluse în structura tezei. Contribuția noastră principală se referă la mineritul diferitelor tipuri de dependențe de date (în date relaționale și XML) utilizând analiza formală a conceptelor; ne concentrăm, de asemenea, pe managementul dependențelor, prevenirea incoerențelor și extragerea de date prin prezentarea modului în care abordările noastre contribuie la descoperirea de cunoștințe conceptuale. În acest rezumat numerotarea capitolelor urmează cea a tezei.

În capitolul 1 introducem principalele contribuții și modul în care acestea sunt structurate în teză.

În capitolul 2 prezentăm fundamentele teoretice, în capitolul 3 prezentăm modul în care FCA poate fi utilizată pentru a caracteriza dependențe funcționale (FD-uri). O primă problemă în utilizarea de FCA pentru a obține dependențele funcționale este faptul că seturile de date sunt în general multivaloarea și nu binare. Acest lucru înseamnă că setul de date trebuie transformat într-un fel pentru obținerea unui context binar ale cărui implicații sunt echivalente cu aceste dependențe. Din păcate acest lucru duce la o nouă reprezentare a datelor, una pătratică în numărul obiectelor referitoare la datele originale, care nu permit aplicarea acestei metode pe seturi de date de mari dimensiuni. Acesta a fost motivația noastră de a crea o metodă de construire a unor fișiere index sau inversate pentru a optimiza construirea contextului formal al dependențelor funcționale. Metoda noastră permite o caracterizare echivalentă dar cu proprietăți de calcul mai bune. În acest capitol prezentăm software-ul nostru (FCAFuncDepMine) care poate fi utilizat în proiectarea bazelor de date relaționale și în detectarea dependențelor funcționare în tabele existente. Acest software poate fi considerat un mediu de asistență a descoperirii cunoștințelor conceptuale, unul care promovează o abordare centrată pe om a proceselor de descoperire a dependențelor și de reprezentare a cunoștințelor obținute.

În capitolul 4 propunem un cadru în care documentul XML este analizat, iar contextul formal care corespunde reprezentării plane a datelor XML este construit. Apoi este analizat

din punctul de vedere al descoperirii cunoștințelor. Latticea conceptuală astfel obținută este o reprezentare grafică utilă a elementelor incluse în documentul XML analizat și a ierarhiei lor. Tot în acest capitol este introdus software-ul nostru (FCAMineXFD). Acest software găsește cheile și dependențele funcționale în datele XML, fiind implicații atribuite în contextul formal construit. Schema documentului XML este transformat în GTT-XNF prin utilizarea dependențelor funcționale identificate.

În capitolul 5 introducem abordarea noastră referitoare la managementul dependențelor și prevenirea contradicțiilor. Descriem strategia noastră de a găsi dependențe funcționale condiționate (Conditional Functional Dependencies - CFD-uri) [30] și reguli de asociere (Association Rules – AR-uri) [2], și în loc de a le utiliza pentru a curăța “date murdare” (dirty data) le folosim pentru a preveni apariția lor în baza de date. Acest lucru a fost realizat prin diferențierea dependențelor stricte de cele aparente. Dacă o dependență este strictă, se iau în considerare toate posibilele excepții de la dependența invalidă, prin urmare la momentul validării sunt generate condiții care previn inserții și actualizări (update) care nu sunt conforme cu dependența validată. Numărul modelelor poate fi foarte mare, de aceea am introdus o relație de ordine parțială “mai utilă” (More Useful - MU), și în fiecare caz de validare sunt luate în considerare supremurile, adică dependențele cele mai utile. Relația de ordine parțială a MU este extrem de utilă deoarece previne inundarea bazei de date cu dependențe valide care sunt consecințele ale unor dependențe mai utile. Abordarea de a preveni apariția de date murdare în baza de date ajută în menținerea consistenței datelor stocate. Împreună cu managementul complet al dependențelor, aplicația noastră implementată, DependencyManager, utilizează metode de Analiza Formală a Conceptelor pentru a analiza dependențele stricte, și pentru a trage concluzii importante, astfel ajutând utilizatorii în prevenirea neconcordanțelor, în repararea erorilor și optimizarea interogărilor și aplicațiilor prin furnizarea unei lattice de dependențe, utilizând utilitatea ca și relație.

În capitolul 6 prezentăm instrumentul nostru care ajută extragerea datelor din diferite tipuri de surse de date. Extractorii de date sunt în general greu de întreținut pentru că trebuie să facă față unor structuri de date posibil infinite și a unor moduri posibil infinite pentru a descrie același lucru. Pentru a simplifica sarcina dificilă de a integra noi surse de date și de a lucra cu schimbări structurale în sursele de date existente am implementat o nouă abordare denumită *RK*. În *RK* arborele semantic generat de regulile semantice (utilizate ca și date de intrare) descrie un model care reprezintă structura sursei de date. Acest instrument este fezabil așa cum este, dar potențialul maxim este atins când se creează o colaborare între proprietarul sursei de date și extractorul, regulile semantice putând fi actualizate ori de câte ori structura sursei de date este schimbată. Acest lucru ar face instrumentul imun la schimbările structurale. Facilitarea extragerii semantice a datelor ar facilita și căutările semantice pe baza datelor extrase. Rezultatele experimentului în care s-a folosit instrumentul *RK* au dovedit fezabilitatea abordării noastre prin faptul că extragerea a fost amplificată din punctul de vedere al preciziei și al vitezei.

1.3 Contribuții publicate anterior

Rezultatele prezentate în această teză au fost deja publicate. Metoda noastră de construire a unor fișiere index sau inversate pentru a optimiza construirea contextului formal al dependențelor funcționale a apărut în [73; 44]. Software-ul nostru, FCAFuncDepMine de extragere a dependențelor funcționale care dețin un set de date prin folosirea de FCA poate fi găsit în [48].

Cadrul nostru care oferă o vizualizare interactivă a explorării dependențelor în date XML poate fi găsit în [47; 43; 74]. Managementul dependențelor și prevenirea incoerențelor sunt, de asemenea, prezentate în [46]. Instrumentul nostru care ajută extragerea datelor din diferite tipuri de surse de date a fost publicat în [45].

2

Fundamentele teoretice

În capitolul al doilea al tezei prezentăm fundamentele teoretice: conceptele de baza ale FCA: ordine și latică, context, concept și latică aferentă, cu exemple și figuri aferente. Scalarea conceptuală și dependențele în baze de date relaționale: funcționale, condițional funcționale și reguli de asociere, dependențele în XML, sunt redate prin definițiile corespunzătoare, exemple și figuri.

3

Caracterizarea dependențelor funcționale (FD) cu contexte formale

În acest capitol prezentăm prima noastră contribuție care se referă la dezvoltarea abordărilor optimizate de extragere FD cu utilizare FCA. Rezultatele prezentate în acest capitol au fost publicate în [73; 44; 48].

O primă problemă în utilizare FCA pentru a calcula FD este faptul că seturile de date sunt în general multivaloarea și nu binare. Acest lucru înseamnă că setul de date trebuie transformat într-un fel pentru obținerea unui context binar ale cărui implicații sunt echivalente cu aceste dependențe. Hereth în [38] oferă o soluție pentru această binarizare. Din păcate în cazul metodei descrise în [18] numărul obiectelor contextului rezultat este pătratic referitor la original, astfel această metodă nu poate fi aplicată pe seturi de date de mari dimensiuni. Aceasta a fost motivația noastră de a crea o metodă de construire a unor fișiere index sau inversate pentru a optimiza construirea contextului formal al dependențelor funcționale. Ne-am propus să descoperim cunoștințe conceptuale, totuși scopul principal a fost de a reduce la minim datele de intrare înainte de a aplica FCA. Am optimizat metoda lui Hereth prin reducerea semnificativă a fișierului de context.

În continuare prezentăm metoda noastră de a construi contextul dependențelor funcționale al unei baze de date tabel.

În procedura de creare a contextului excludem perechile de tuple care sunt simetrice și reflexive pentru a evita redundanțele. În partea superioară a lăței conceptuale se vor afla perechi de tuple în care nu sunt valori comune pentru atributele corespunzătoare, astfel putem omite generarea acestor perechi de tuple. Perechile cu forma $(t; t)$, unde t este un

tuplu al tabelului, au toate atributele în comun. Aceste obiecte vor ajunge în partea de jos a lăței, astfel putând fi omise, fiindcă nu schimbă aceste implicații în context. Pentru a optimiza construirea contextului formal al dependențelor funcționale am construit fișiere index sau inversate pentru valorile fiecărui atribut. Cu aceste indici inversate numărul rândurilor în fișierul de date al contextului formal care rezultă în urma aplicării metodei noastre este jumătatea valorii care rezultă în urma aplicării metodei lui Hereth în [38]. În consecință putem reduce timpul necesar pentru a construi lățea conceptuală pentru dependențe funcționale și putem elimina dependențele inutile.

3.1 Rezultatele experimentului

Pentru a evalua metoda noastră de creare a contextului prin folosirea indicilor inversate am efectuat experimente pe diferite seturi de date.

Se consideră următoarea tabelă incorect concepută:

```
StudMarks [StudID, StudName, GroupID, Email, SpecID,  
            SpecName, Language, DiscID, DName, CreditNr, Mark]
```

Această tabelă conține 19500 de rânduri și 11 coloane; metoda noastră a redus numărul obiectelor de la 380.250.000 la 110.073.342 în fișierul de context.

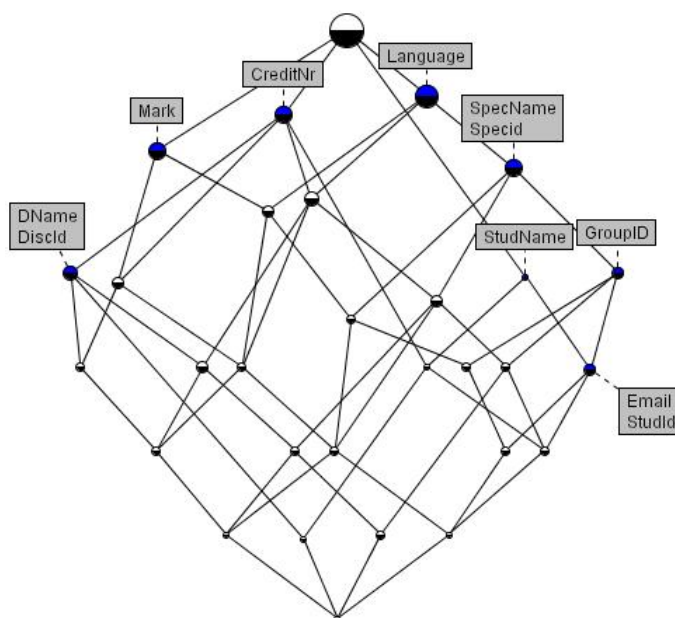


Fig. 3.1: Lățea conceptuală $FD(StudMarks, \vec{\mathbb{K}}(Uni))$

Am aplicat FCA pe contextul redus; Figura 3.1 include lățea conceptuală rezultată în cazul tabelii StudMarks. Implicațiile în această lățea, care sunt dependențe funcționale în tabelă, pot fi interpretate în următorul fel: conceptul cu denumirea *GroupID* este un

subconcept al conceptului cu denumirea *SpecID* și *SpecName*. Acest lucru înseamnă că în fiecare pereche de tuple, unde câmpul *GroupID* are aceeași valoare, specializarea ID și numele este identic. Astfel avem următoarele implicații, care sunt dependențe funcționale:

$$GroupID \rightarrow SpecID, SpecName$$

În același mod următoarele dependențe funcționale pot fi citite din laticea conceptuală:

$$DiscID \rightarrow CreditNr$$

$$DName \rightarrow CreditNr$$

$$SpecID \rightarrow Language$$

$$SpecName \rightarrow Language$$

$$StudID \rightarrow StudName$$

$$Email \rightarrow StudName$$

De asemenea am efectuat experimente utilizând bine-cunoscutele seturi de date din baza de învățare automată UCI <http://archive.ics.uci.edu/ml/>, și anume seturile de date Wisconsin Breast Cancer (WBC) și Mushroom. Putem concluziona că prin utilizarea acestor indici inversate, numărul rândurilor în fișierul de date al contextului formal care rezultă în urma aplicării metodei noastre este jumătatea valorii care rezultă în urma aplicării metodei lui Hereth în [38]. Tabelul 3.1 include sumarul performanței metodei noastre. Astfel putem reduce timpul necesar pentru a construi laticea conceptuală pentru dependențe funcționale și putem elimina dependențele inutile. La acest punct subliniem faptul că prin optimizarea dimensiunii contextului calitatea setului de dependențe creat nu s-a deteriorat; metoda noastră permite o caracterizare echivalentă dar cu proprietăți de calcul mai bune. Am implementat metoda prezentată în acest capitol și l-am completat cu un instrument de software care analizează tabelele de date relaționale existente.

Dataset	Arity	No. Rows	Nr. obiectelor cu metoda noastră	Nr. obiectelor cu metoda lui Hereth [38]
Mushroom	22	8,124	32,995,626	65,999,376
WBC	11	699	210,934	481,636
StudentMarks	11	19,500	110,073,342	380,250,000

Table 3.1: Rezultate experimentale: contextul optimizat

Software-ul nostru denumit FCAFuncDepMine, construiește un context formal optimizat de dependențe funcționale. Software-ul poate fi utilizat în proiectarea bazelor de date relaționale și în detectarea dependențelor funcționare în tabele existente. În prezent FCAFuncDepMine suportă trei tipuri de format de fișiere de date de intrare: XML, format csv și tabele de baze de date relaționale. FCAFuncDepMine salvează contextul formal în format .cex utilizat pe scară largă.

Optimizările au loc în FCAFuncDepMine în mod adecvat pentru a transforma seturi de date de mari dimensiuni în contexte formale. Software-ul nostru poate fi considerat un mediu de asistență a descoperirii de cunoștințe conceptuale, care promovează o abordare centrată pe om a proceselor de descoperire, a dependențelor și de reprezentare a rezultatelor.

4

Dependențe funcționale XML cu FCA

Ca o a doua contribuție propunem un cadru în care este analizat documentul XML, și este construit contextul formal, care corespunde reprezentării plane a datelor XML. Rezultatele prezentate în acest capitol au fost publicate în [47; 43; 74]. Metoda noastră este suportată de un cadru denumit FCAMineXFD care găsește cheile și dependențele funcționale în datele XML, fiind implicații atribuite în contextul formal construit. Latticea conceptuală obținută este o reprezentare grafică utilă a elementelor incluse în documentul XML analizat și a ierarhiei lor. Schema documentului XML este transformată în GTT-XNF [90] prin utilizarea dependențelor funcționale identificate.

Pentru a realiza acest lucru datele XML trebuie convertite în contexte formale FCA. Software-ul nostru poate analiza întregul document XML sau o clasă tuplu C_p dat de traiectoria p . Datele XML pot fi convertite într-o relație total neimersată, o singură tabelă relațională și algoritmi de descoperire FD existenți pot fi aplicați în mod direct. Dintr-un document XML care conține schema datelor de la început, creăm arbori de tupluri generalizate.

Ca un prim pas trebuie să definim obiectele și atributele de interes și să creăm modele de XML în termeni de context.

- Alegerea atributelor FCA: *PathEnd/ElementName*
 - în cazul nodurilor interne (non-leaf) numele atributului este construit în următorul fel:
<ElementName>+”@key” ” și valoarea lui va fi valoare cheii asociate. Mai multe

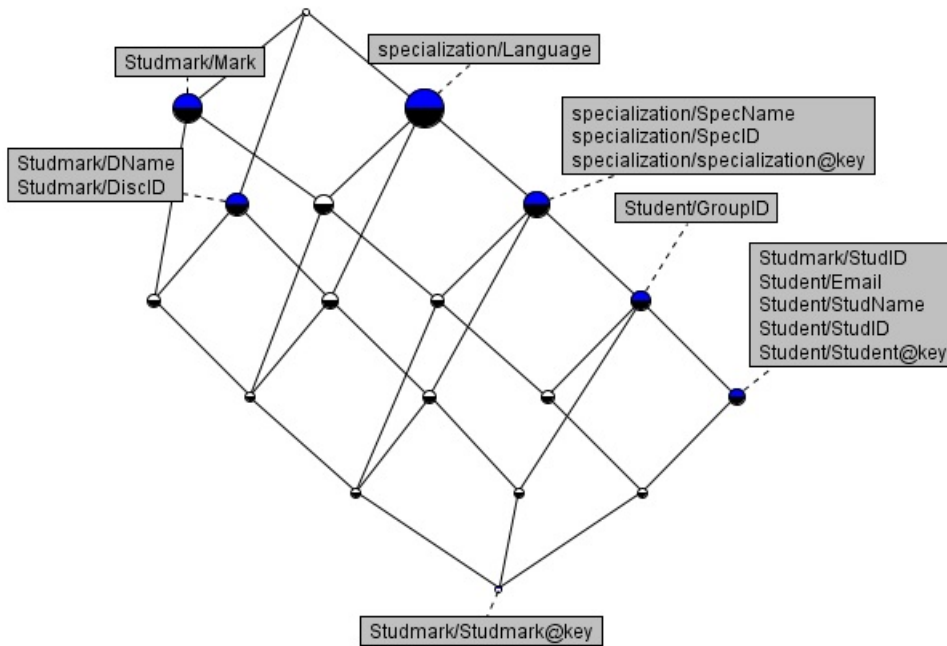


Fig. 4.1: Latticea conceptuală pentru $C_{specialization}$

elemente care au aceeași traiectorie, vor avea același nume de atribut, dar valorile vor fi diferite.

– frunzele (nu valorile pe frunze, ci numele elementelor pe frunze) arborelui tuplu.

- Alegerea obiectelor: obiectele sunt considerate perechile tuplu ale arborelui, mai precis perechile de tuplu din tabela plată. Valorile cheie asociate elementelor interne (non-leaf) și valorile elementelor leaf sunt utilizate în aceste perechi de tupluri.

Documentul XML analizat poate avea un număr mare de arbori tuplu. Folosim aceleași tehnici de optimizare ca și în capitolul anterior, filtrăm perechile de tupluri și oțitem acele perechi care nu au atribute comune, ceea ce nu modifică ierarhia conceptuală.

Odată ce obiectele și atributele contextului sunt definite, generăm conceptele și creăm latticea conceptuală în același mod ca și în cazul bazelor de date relaționale. O lattice conceptuală pentru o bază de date XML universitară (clasa tuplu $C_{specialization}$) se poate observa în figura 4.1, unde putem interpreta fiecare concept și genera lista tuturor dependențelor funcționale.

Example 4.1. În nodul conceptual cu eticheta $specialization/specialization@key$, $specialization/SpecID$, $specialization/SpecName$ obiectele asociate sunt perechile de tupluri, unde valorile pentru $specialization=SpecID$ sunt aceleași. Deci avem următoarele FD-uri XML:

$$\langle C_{specialization}, ./SpecID, ./SpecName \rangle$$

$$\langle C_{specialization}, ./SpecID, ./specialization@key \rangle$$

$$\langle C_{specialization}, ./specialization@key, ./SpecID \rangle$$

$$\langle C_{specialization}, ./specialization@key, ./SpecName \rangle$$

$$\langle C_{specialization}, ./SpecName, ./specialization@key \rangle$$

$$\langle C_{specialization}, ./SpecName, ./SpecID \rangle$$

Din laticea conceptuală se poate citi relația dintre concepte. Există 1:N relații de la *Specialization* la *Group*, de la *Group* la *Students*, de la *Students* la *Studmark*. *Disciplines* este într-o relație N:N cu *Students*, conectat la nodul *Studmark* în acest caz. În abordarea noastră pe lângă XFD-uri putem găsi cheile XML pentru un document XML dat.

Implicațiile găsite de FCAMineXFD conțin unele FD-uri cu RHS ca și valori cheie (./@key). Acestea pot fi utilizate pentru a detecta cheile în XML. Acesta este un pas important în descoperirea de cunoștințe bazate pe FCA, deoarece până acum au putut fi citite doar dependențele din laticea conceptuală, dar prin utilizarea metodei noastre putem găsi și cheile.

Sunt de exemplu două FD-uri (Ex. 4.1) cu RHS ca și ./specialization@key în clasa tuplu $C_{specialization}$, astfel cheile XML detectate sunt:

$$\langle C_{specialization}, ./SpecID \rangle, \langle C_{specialization}, ./SpecName \rangle.$$

Dat fiind setul de chei și setul de dependențe descoperite de instrumentul nostru, vom converti schema XML într-una corectă.

Ca o a doua contribuție propunem un cadru în care este analizat documentul XML, este construit contextul formal care corespunde reprezentării plane a datelor XML și este efectuată analiza acestuia din punctul de vedere al descoperirii de cunoștințe. Așa cum am observat laticea conceptuală obținută în așa fel este o reprezentare grafică valoroasă a elementelor incluse în documentul XML analizat, și a ierarhiei lor. Software-ul nostru, FCAMineXFD, găsește cheile și dependențele funcționale între datele XML, fiind implicații atribuite în contextul formal construit. Schema documentului XML este transformat în GTT-XNF prin utilizarea dependențelor funcționale identificate.

5

Managementul dependențelor și prevenirea incoerențelor

În acest capitol descriem a treia contribuție care găsește dependențe funcționale condiționate (Conditional Functional Dependencies - CFDs) [30] și reguli de asociere (Association Rules - ARs) [2] și în loc de a le utiliza pentru a curăța “date murdare” (dirty data) le folosim pentru a preveni apariția acestora în baza de date. Acest lucru a fost realizat prin diferențierea dependențelor stricte de cele aparente. Dacă știm despre o dependență că va deveni validă în viitor, ne putem baza pe acestea prin crearea unor limite care garantează că regula CFD nu va fi încălcată de inserții sau modificări. Împreună cu managementul complet al dependențelor aplicația noastră implementată denumită DependencyManager, utilizează metode de Analiza Formală a Conceptelor pentru a analiza dependențele stricte și trage concluzii importante, astfel ajutând utilizatorii în prevenirea neconcordanțelor, în repararea erorilor și optimizarea interogărilor și aplicațiilor prin furnizarea unei lăți de dependențe, folosind utilitatea ca și relație. Managementul dependențelor și prevenirea incoerențelor sunt, de asemenea, prezentate în [46].

Aplicația noastră învață FD-uri, CFD-uri și AR-uri dintr-o bază de date arbitrară. Putem valida dependențele învățate prin acceptarea sau respingerea lor, astfel modificând datele de la candidat la dependențe stricte (SD) sau aparente (AD). Am creat sprijin pentru simplificarea validării prin utilizarea relației MU astfel încât cei mai utili candidați să fie luați în considerare dintre toate celelalte dependențe, de aceea reducem în mod drastic numărul dependențelor care trebuie luate în considerare și păstrăm doar cele mai utile SD-uri deoarece toate celelalte dependențe sunt implicate de acestea. Rezultatul cercetării noastre este un

sistem care este capabil de a preveni toate incoerențele posibile care rezultă din inserții sau modificări (update) care nu sunt conforme schemei SD-urilor acceptate prin utilizarea limitelor. Folosim metode FCA pentru a analiza CFD-urile/AR-urile stricte și pentru a trage concluzii utile. Cu alte cuvinte: exploatăm FCA-ul pentru a obține concepte care conțin ontologie.

Strategia noastră se bazează pe învățarea automată, validarea semi-automată și uitarea automată. Am definit un algoritm care este folosit de aplicație pentru a învăța noi CFD-uri și AR-uri. Acceptarea automată a tuturor dependențelor învățate poate fi riscantă mai ales dacă datele sunt critice. Validarea nu este întotdeauna ușoară pentru utilizatori, deoarece pot fi multe dependențe învățate, dar am utilizat relația “mai util” (MU), care reduce semnificativ dificultatea validării dacă este folosită în mod corespunzător. Așa cum am mai menționat, utilizatorul vede o listă cu cele mai utile CFD-uri/AR-uri și poate decide dacă acestea sunt SD-uri sau AD-uri. Dacă utilizatorul acceptă o dependență D_1 atunci toate D_2 vor deveni instantaneu invalizi ceea ce îndeplinește criteriile de $D_1 MU D_2$. Ca urmare, orice D_2 mai puțin util decât D_1 este redundant cu D_1 . Dacă utilizatorul respinge un CFD/AR D_1 , atunci acesta va deveni invalid iar toate CFD-urile/AR-urile D_2 care îndeplinesc criteriile

$$(D_1 MU D_2) \wedge (\nexists D_3 \text{ astfel ca } D_1 MU D_3 MU D_2)$$

vor fi afișate utilizatorului deoarece vor deveni cele mai utile AR-uri/CFD-uri.

Am estimat beneficiile de utilizare a relației MU . Dacă o tabelă are n coloane și o dependență are o coloană pentru condiția, s coloane în setul de coloane determinante, și d coloane în setul de coloane dependente, există $2^{(n-s-d-1)} + 2^d - 1$ posibile dependențe care sunt mai puțin utile. Dacă acceptăm dependența dată, în mod automat refuzăm toate dependențele mai puțin utile, ceea ce optimizează, ajutând sistemul/utilizatorul cu reducerea automată a dependențelor care trebuie luate în considerare cu un maxim de $2^{(n-s-d-1)} + 2^d - 1$ dependențe. Această metodă ajută mult în validarea și în utilizarea dependențelor. Aceasta este o optimizare exponențială.

5.1 Rezultatele experimentului

Am evaluat eficiența și eficacitatea algoritmului nostru pe patru seturi de date. În seturile de date denumite Numbers1 și Numbers2 rândurile au fost numere generate folosind formule pentru a garanta apariția CFD-urilor și AR-urilor. De asemenea am efectuat experimente utilizând seturi de date reale din baza de învățare automată UCI <http://archive.ics.uci.edu/ml/>, și anume seturile de date Wisconsin Breast Cancer (WBC). Setul de date Exceptions a fost, de asemenea, generat de noi, prin folosirea datelor reale.

Tabelul 5.1 descrie parametrii seturilor de date utilizate, rata MinOccurrenceFrequency și MinOccurrenceRate utilizată în experimentul nostru și numărul de dependențe accep-

Dataset	Arity	No. Rows	MOF	MOR	SD	AD
Numbers1	6	20000	200	0.15	95	13
Numbers2	5	200000	650	0.2	118	11
WBC	11	699	37	0.06	242	21
Exceptions	5	11192	1000	0.1	43	8

Table 5.1: Rezultatele experimentului (MOF = MinOccurrenceFrequency, MOR = MinOccurrenceRate)

tate/respinse: SD-uri și AD-uri pentru fiecare set de date.

Pe baza analizei noastre am reușit să reducem semnificativ numărul dependențelor luate în considerare, am constatat că doar 10% -30% dintre dependențe au fost utile.

În cazul în care un model SD este acceptat (validat), o limitare împiedică inserțiile și actualizările incoerente cu modelul acceptat. Am cuantificat cantitatea de date protejate de sistemul nostru împotriva încălcării regulii de incoerență. În setul nostru de date Numbers1 8620, Numbers2 102034, BCW 384, Exceptions 7083 reprezintă numărul de înregistrări care îndeplinesc cel puțin un SD, cu o inserție sau o actualizare arbitrară în viitor având o probabilitate de 43,1%, 51,01%, 54,93% și 63,28% de a fi protejat prin sistemul nostru împotriva incoerențelor care încalcă SD-urile noastre. Dacă apare o eroare din cauza manipulării necorespunzătoare a datelor unde un SD este aplicabil, problema poate fi identificată în jurnalele de eroare și, în mod potențial, erorile pot fi găsite și reparate.

În DependencyManager am implementat un modul FCA, care studiază proprietățile SD-urilor descoperite cu ajutorul aplicației noastre și trage concluzii interesante și cunoștințe adiționale despre ele. Abordarea noastră a fost agnostică față de natura conținutului bazei de date; am implementat un FCA aplicabil dependențelor acceptate ale fiecărei baze de date prin utilizarea sistemului nostru. Ca un prim pas am transformat cunoștințele existente într-un context FCA bazat pe SD-uri. În acest context obiectele sunt FD-uri, iar atributele sunt proprietățile posibile ale obiectelor. Atributele sunt Weakness (Slăbiciune), Determinant Column Cluster's Size (dimensiunea fasciculului coloanei determinante), Column Cluster Size (dimensiunea fasciculului coloanei), Frequency of Occurrence (frecvența apariției) și Almost Symmetrical (aproape simetric). Am determinat seturile fuzzy de None (niciunul), Very Low (foarte scăzut), Low (scăzut), Medium (mediu), High (mare), Very High (foarte mare) și Total (total) pentru acest atribut. În al doilea rând am creat conceptele (nu doar prin gruparea obiectelor care au același set de atribute, dar și prin compararea valorilor lor), și apoi le-am transformat într-o latice conceptuală (fără binarizare), ceea ce a fost baza analizelor ulterioare. Cu alte cuvinte: exploatăm FCA-ul pentru a obține concepte care conțin ontologie. FCA este relevant în acest studiu prin furnizarea informațiilor conceptuale, ceea ce este rezultatul analizei SD cu ajutorul atributurilor prezentate anterior.

După cunoștința noastră, niciun studiu anterior nu a folosit implementarea FCA, cu obiectele de intrare fiind SD-uri și atributele de intrare fiind proprietățile acestor dependențe.

Pe lângă managementul complet al dependențelor strategia noastră previne apariția

incoerențelor în locul curățării celor existente. Putem concluziona că aplicația noastră este accesibilă, ușor de utilizat și eficace în prevenirea incoerențelor. Algoritmul nostru reprezintă o abordare nouă în managementului CFD și în prelucrarea FCA, prin analiza modelelor de baze de date abstracte.

6

Extragerea datelor prin utilizarea arborilor semantici

Un rezultat major al prezentei cercetări este un instrument rapid și precis care ajută extragerea datelor semantici din diferite tipuri de surse de date. Pentru a simplifica sarcina dificilă de a integra noi surse de date și de a lucra cu schimbări structurale în sursele de date existente, am implementat o nouă abordare numită *RK*. *RK* este un instrument care funcționează ca un ghid pentru scopuri de extracție; arborele semantic generat de regulile semantice (utilizate ca date de intrare) descrie un model care reprezintă structura sursei de date.

Scopul cercetării noastre este de a simplifica procesul complicat de creare a unui extractor de date. Ne propunem să motivăm proprietarii de surse de date să ajute echipele de extragere a datelor în efectuarea sarcinilor de extragere a datelor din sursele de date. Pentru a atinge acest obiectiv am creat *RK*, un instrument care utilizează un limbaj interpretat semantic pentru a defini regulile semantice care reprezintă conceptele, inclusiv traiectoria lor structurală, relațiile de tip părinte-copil, pluralitatea, setul de sinonime/cuvinte cheie, prioritate și alte atribute care descriu conceptul.

Utilizăm termenul de concept aici pentru a desemna toate regulile cu același cuvânt cheie. De exemplu în cazul în care informațiile care trebuie extrase se referă la sectorul imobiliar, posibilele concepte pot include mărimea sau prețul. Din aceste reguli *RK* generează obiecte care servesc ca noduri în construirea arborelui semantic. Arborele descrie modelul structural-semantic utilizat ca un ghid în extragerea funcționalităților. Aceste metode, împreună cu funcționalitățile de navigare sunt definite în scopuri generale, totuși pot fi personalizate,

iar în cazul funcțiilor de navigare pot fi suprascrise dacă acest lucru este necesar. Sistemul nostru suportă extragerea datelor din surse de date structurate ierarhic, cum ar fi pagini web printre altele.

Principalele beneficii ale propunerii noastre includ potențialul de a ajuta căutările conceptuale și posibilitatea unor schimbări structurale fără necesitatea de a schimba aplicațiile de extragere a datelor care extrag din surse de date ale căror structură a fost modificată.

Instrumentul nostru susține portabilitatea, personalizarea și extragerea conceptuală, semantică. Implementarea actuală poate fi integrată în paginile web și are scopul de a da instrucțiuni crawler-ilor. După verificarea pentru a vedea dacă *RK* este integrat în sursa de date, extractorul poate extrage întregul set de date prin aplicarea unei singure funcțiuni. Pentru a sprijini definirea regulilor semantice am dezvoltat limbajul semantic *RK*.

Fiecare regulă reprezintă un concept și este conștientă de conceptul părinte sau predecesor, poziția sa în raport cu poziția structurală a conceptului său părinte sau succesori și poziția structurală în raport cu valoarea conceptului. Polimorfismul este suportat, fiindcă același tip de concepte pot avea mai multe forme posibile în structură. Acest lucru înseamnă că mai multe reguli pot descrie același concept pentru a gestiona diferitele cazuri ale apariției sale. În astfel de cazuri acestea sunt identificate cu ajutorul indicilor. Un concept poate avea mai mulți părinți sau predecesori potențiali. Sinonimele sunt susținute din cauza divergențelor culturale. Nu există nicio ontologie standard; denumirea conceptelor diferă de la cultură la cultură și multe concepte au diferite denumiri în aceeași cultură. Instrumentul respectă diversitatea culturală, oferind posibilitatea de a defini sinonime, adică diferite cuvinte cheie pentru același concept.

Regulile sunt salvate în fișiere *.rk*. Personalizarea extragerii este posibilă prin combinarea acestora. Numim aceste combinații un tip de căutare. În mod ideal aceste reguli sunt create și întreținute de către proprietarii surselor de date, dar pot fi definite de către oricine care are acces la fișierele *.rk* sau pot fi generate în mod automat de către un modul. Extragerea funcționează în baza regulilor structural-semantică. Dacă reclamele nu sunt considerate a fi un conținut relevant, nicio regulă semantică nu va fi definită pentru a le suporta, astfel extractorii nici nu vor lua în considerare aceste părți.

Regulile sunt convertite în obiecte. Relația de tip părinte-copil, predecesor-succesor a regulilor presupune o ierarhie. Prin crearea unui nod rădăcină abstract, care are toate conceptele cu predecesori nespecificați ca și copii, un arbore este definit, unde toate nodurile corespund regulilor și nodurile sunt relația de predecesor-succesor. Arborele semantic descrie planul de lucru structural-semantic pentru extractor. Modificările structurale au fost întotdeauna greu de gestionat cu ajutorul extractorilor. În cazul în care se folosește instrumentul nostru, în cazul modificărilor structurale doar regulile semantice trebuie păstrate. Extractorii care folosesc sursa de date nu vor simți efectele acestei schimbări (ceea ce altfel ar face). În schimb ei pot extrage data ca și cum nu ar exista modificări structurale.

După cum am văzut, arborele semantic descrie structura datelor. Datorită structurii ierarhice a arborelui semantic, extractorul poate restrânge în mod repetat spațiul de căutare,

fiindcă din punct de vedere structural un concept include sub-conceptele sale, astfel când se caută sub-conceptele unui concept, reprezentarea structurală a conceptului poate servi ca o sursă a căutării.

Sarcina de extragere este îndeplinită prin simplificarea repetată a spațiului de căutare. În cazul în care un concept nu este plural, în prima fază este găsită în interiorului părintelui sau antecedentului, iar alte cazuri nu vor fi căutate. Dacă este unul plural, se efectuează o căutare completă pentru a identifica toate aparițiile. Cu ajutorul instrumentului nostru extragerea nu este mai complicată decât un apel de funcție. Instrumentul poate fi integrat într-o sursă de date prin copierea codului sursă și prin crearea și publicarea fișierelor .rk. În unele cazuri funcțiile de navigare trebuie suprascrise.

6.1 Rezultatele experimentului

Am efectuat experimente pe diferite site-uri web reale. După verificarea pentru a vedea dacă *RK* este integrat în sursa de date, extractorul a extras întregul set de date prin aplicarea unei singure funcțiuni. O funcție suplimentară a fost utilizată pentru a trece la pagina următoare. Tabelul 6.1 include sumarul rezultatelor experimentului. Toate experimentele au fost efectuate pe un calculator echipat cu un Pentium Dual-Core CPU T4300, 2,1 GHz, cu 4 GB RAM.

	No. of rules	No. of Extracted Objects	Total Download Time (millisec)	No. of Objects/Page	Total Extraction Time (millisec)	Extraction Time/Page (millisec)	Extraction Time/Object (millisec)
Test RealEstate grid page	17	11376	513748	10	38372	33.72	3.37
Test RealEstate details page	21	11376	7566787	10	265288	233.12	23.32
www.boatsandoutboards.co.uk	6	7721	2586972	20	55070	148.84	7.13
http://prep-adm.cloudapp.net/	13	160012	6840503	20	521756	65.13	3.25

Table 6.1: Rezultatele experimentului

Unul dintre site-urile web testate a fost *www.boatsandoutboards.co.uk*, din care am extras *Image* (image), *Link* (link), *Title* (titlu), *Properties* (proprietăți), și *Description* (descriere) pentru bărci. În acest caz sursa de date a conținut 7721 bărci sau accesorii, iar timpul mediu de extragere/obiect a fost de 7,13 milisecunde. În timpul procesului de extragere instrumentul nu efectuează operațiuni care nu sunt necesare, astfel viteza de extragere este optimă. În loc de a evalua sute sau chiar mii de expresii obișnuite (ca și în [92]), în întreaga pagină web se evaluează doar câteva zeci de reguli structural-semantice, se alcătuieste un arbore semantic cu ajutorul căruia toate (și exclusiv) datele necesare sunt extrase.

Un alt site pe care am testat instrumentul nostru de extragere a fost *http://prep-adm.cloudapp.net*, care conține rezultatele examenelor finale ale tuturor absolvenți ai ciclului primar din România. Din această pagină am extras cu succes *County* (județ), *School* (școala), *Name* (numele) și notele a 160012 elevi. Pentru acest experiment am definit 13

reguli structural-semantice. Timpul mediu de extragere pe pagină a fost de 65,13 milise-cunde, ceea ce corespunde la 3,25 milise-cunde/elev. Acest site oferă posibilitatea de a căuta pe ea, dar dacă cineva dorește să analizeze întregul sistem de învățământ din România, are nevoie de toate datele de pe acest site, ceea ce nu este accesibil. Metoda noastră oferă o soluție prin extragerea semantică a tuturor datelor necesare, ceea ce permite aplicarea unei analize suplimentare a datelor.

Nu se raportează precizie și rechemare calculată față de calitatea datelor extrase deoarece nu s-au înregistrat pierderi în experimentele noastre. Precizia extragerii a fost de 100% atâta timp cât regulile structural-semantice au fost definite în mod corect. Timpul de extragere a fost egal cu timpul cumulat necesar pentru crearea arborelui semantic plus extragerea tuturor datelor prescrise cu ajutorul jQuery, mai puțin de 10 milise-cunde, în medie (pe un calculator simplu cu parametrii descrise anterior). Desigur am putea utiliza un super-server și multithreading pentru a obține rezultate și mai bune.

Putem concluziona că instrumentul pe care l-am creat simplifică extragerea datelor, face posibilă crearea mai multor strategii de extragere și facilitează căutări semantice în sistemele care utilizează instrumentul. Extragerea este optimă deoarece cu ajutorul unor reguli se-mantice corecte, algoritmul restrânge în mod repetat spațiul de căutare pentru a crește performanța căutării. În locul analizei textului am oferit o alternativă căutării ierarhice care funcționează în mod constant în subarborele unui anumit concept pentru a găsi valorile sub-conceptelor în baza arborelui semantic alcătuit din seturile de reguli. Cu aplicarea regulilor corecte, datele extrase trebuie să fie corecte și dacă acest lucru este realizat într-un mediu de colaborare, precizia poate fi menținută pe termen lung.

Rezultatele experimentului în care s-a folosit instrumentul RK au dovedit fezabilitatea abordării noastre prin faptul că extragerea a fost amplificată din punctul de vedere al preciziei și al vitezei.

7

Concluzii și planuri de viitor

În această teză am prezentat modul în care FCA-ul poate fi utilizat pentru a caracteriza dependențe funcționale și pentru a extrage cunoștințe din baze de date. În viitor ne propunem să caracterizăm dependențele funcționale utilizând formalismul structurilor model [10], o extensie a FCA-ului clasic pentru a gestiona date complexe. Ne propunem de asemenea să dezvoltăm un software care va crea un tricontext al unui arbore XML și va oferi dependențele funcționale din arborele XML. Pe de altă parte am dori să extragem, în ceea ce privește FCA, și alte tipuri diferite de dependențe care pot include un anumit set de date (nu numai în cele relaționale și XML, dar și în baze de date NoSQL).

În capitolul 5 am introdus DependencyManager, abordarea noastră referitoare la managementul dependențelor și prevenirea contradicțiilor. În viitor ne propunem să facem această strategie și mai utilă, ne propunem să găsim și să analizăm candidați SD pe tabele ordonate transversal. Dorim de asemenea să generalizăm un set de condiții prin utilizarea mai multor coloane în funcțiile booleene în locul unei singure și cu utilizarea mai multor tipuri de operatori (implementarea noastră curentă consideră egalitatea ca singurul operator condițional). Automatizarea procesului de validare ar fi utilă; această funcție va fi probabil bazată pe reguli determinate de utilizator care ajută modulul care învață să determine în mod automat dacă un candidat SD este un AD sau un SD. Sistemul salvează erorile care provin din inserții și actualizări (update) eșuate care au fost împiedicate de limitările create pentru SD-uri; acesta va genera cunoștințe care pot fi utilizate să determine cauza eșecului inserțiilor și actualizărilor. În cazul în care eșecul este un model incorect acceptat, modelul trebuie respins; în caz contrar acțiunea de utilizator eronată sau funcționalitatea incorectă poate fi detectată. SD-urile cunoscute pot fi utilizate în creșterea performanței aplicației.

Ca urmare spațiul de stocare/performața de siguranță (backup) poate fi optimizat prin crearea modelelor/șabloanelor în loc de a stoca toate cazurile de date redundante care corespund modelelor valide. Modelele valide pot fi arhivate (cache) care permit sistemele-client să deducă valori și astfel să reducă numărul cererilor de la baza de date. Prin urmare, dimensiunea răspunsurilor bazei de date pot fi reduse, optimizând comunicarea între sistemul client, serverul aplicației și serverul bazei de date.

În capitolul 6 am introdus instrumentul nostru *RK* care ajută extragerea datelor din diferite tipuri de surse de date. În prezent setul de reguli în *RK* poate fi definit ușor, dar regulile trebuie introduse manual, deoarece niciun modul de definire automată a regulilor nu a fost încă creat. În viitor ne propunem să creăm accesorii (add-on) de browser și aplicații desktop care suportă definiții de reguli grafice generatoare de reguli RK pe baza acțiunii utilizatorului. Semantic Tree Builder (Creatorul de arbori semantici) citește regulile din setul de reguli și creează un arbore semantic, creând un nod de rădăcină abstract și alte noduri, care formează un arbore. Nodul de rădăcină există pentru a garanta faptul că structura datelor va fi un arbore, iar celelalte reprezintă regulile. Relațiile dintre noduri formează vârfurile arborelui. În prezent nu este nicio garanție referitor la faptul că graficul nu va conține cicluri, ceea ce înseamnă că în cazul în care un dezvoltator creează un ciclu în mod accidental, sistemul nu îl împiedică.

În viitor ne-am propus să creăm un validator pentru arborele semantic, ceea ce ar considera arborele să fie valid doar dacă acesta este într-adevăr un arbore, și nu conține cicluri. Modulul Extractor al *RK* extrage reprezentarea textuală a conceptelor situate în arborele semantic. Dorim să implementăm un layer superior care va folosi Extractorul și va suporta extragerea altor tipuri de meta-date, cum ar fi fișiere sunet, video, imagini sau alte tipuri pe baza datelor textuale extrase cu ajutorul Extractorului. Extractorul nu suportă analiza datelor textuale; layerul superior va suporta și analiza textuală. Ideile din spatele modulelor Semantic Tree Builder, Extractor și Navigator sunt agnostice din punct de vedere tehnologic; suportă extragerea din orice sursă de date având structura ierarhică a datelor (implementarea noastră curentă suportă extragerea de pe pagini web). Ne propunem, de asemenea, suportul altor tipuri de surse de date.

Bibliografie

- [1] S. Abiteboul, R. Hull, V. Vianu. *Foundations of Databases*, Addison-Wesley, 1995
- [2] R. Agrawal, R. Srikant. Fast algorithms for mining association rules in large databases, In *VLDB '94: Proceedings of the 20th International Conference on Very Large Data Bases*. Morgan Kaufmann Publishers Inc., San Francisco, CA, USA, 1994, 487–499
- [3] S. Andrews. In-Close, a Fast Algorithm for Computing Formal Concepts, In Rudolph, Dau, Kuznetsov (Eds.), *Supplementary Proc. of ICCS'09*, CEUR WS 483, 2009
- [4] M. Arenas, L. Libkin. A normal form for XML documents, *TODS* 29(1), 2004, 195–232
- [5] M. Arenas, L. Libkin. An information-theoretic approach to normal forms for relational and XML data, *JACM* 52(2), 2005, 246–283
- [6] M. Arenas, L. Libkin, W. Fan. On the Complexity of Verifying Consistency of XML Specifications, *SIAM J. Comput.* 38(3), 2008, 841–880
- [7] J. Baixeries. A Formal Concept Analysis framework to model functional dependencies. *Mathematical Methods for Learning*, 2004
- [8] J. Baixeries. A formal context for symmetric dependencies, In R. Medina, S. A. Obiedkov, editors, *ICFCA*, volume 4933 of LNCS, Springer, 2008, 90–105
- [9] J. Baixeries. *Lattice Characterization of Armstrong and Symmetric Dependencies* (PhD Thesis), Universitat Politècnica de Catalunya, 2007
- [10] J. Baixeries, M. Kaytoue, A. Napoli. Characterizing Functional Dependencies in Formal Concept Analysis with Pattern Structures, *Annals of Mathematics and Artificial Intelligence*, 2014
- [11] P. Becker, J. Hereth, G. Stumme. ToscanaJ: An open source tool for qualitative data analysis, In V. Duquenne, B. Ganter, M. Liquiere, E. M. Nguifo, and G. Stumme, editors, *Advances in Formal Concept Analysis for Knowledge Discovery in Databases*, 2002.

- [12] R. Belohlavek, M. Trnecka. Basic Level in Formal Concept Analysis: Interesting Concepts and Psychological Ramifications, IJCAI'13 Proc. of the Twenty-Third international joint conference on Artificial Intelligence, 2013, 1233–1239
- [13] G. Beskales, I. F. Ilyas, L. Golab. Sampling the repairs of functional dependency violations under hard constraints, PVLDB, 3(1): 2010, 197–207
- [14] P. Bohannon, W. Fan, F. Geerts, X. Jia, A. Kementsietsidis. Conditional functional dependencies for data cleaning, In ICDE, IEEE, 2007, 746–755
- [15] P. V. Borza, O. Sabou, C. Sacarea. OpenFCA, an open source formal concept analysis toolbox, IEEE International Conference on Automation Quality and Testing Robotics AQTR, Vol 3, 2010, 1–5
- [16] P. Buneman, S. Davidson, W. Fan, C. Hara, W.-C. Tan. Keys for XML, Proc. WWW, Hong Kong, China, 2001, 201–210
- [17] A. Buzmakov, S. O. Kuznetsov, A. Napoli. Scalable Estimates of Concept Stability, ICFCA, 2014, 157–172
- [18] S. Chakravarthy, H. Zhang. Visualization of association rules over relational DBMS, SAC '03: Proceedings of the 2003 ACM Symposium on Applied Computing, New York, NY, USA, 2003, 922–926
- [19] C. H. Chang, M. Kayed, M. R. Girgis, K. F. Shaalan, A survey of web information extraction systems, IEEE Transactions on Knowledge and Data Engineering, 18(10), 2006, 1411–1427
- [20] F. Chiang, R. J. Miller. Discovering data quality rules, PVLDB, 1(1): 2008, 1166–1177
- [21] E. F. Codd. The relational model for database management: version 2. Addison-Wesley Longman Publishing Co., Inc., Boston, MA, USA, 1990
- [22] G. Cong, W. Fan, F. Geerts, X. Jia, S. Ma. Improving data quality: Consistency and accuracy, VLDB, ACM, 2007, 315–326
- [23] G. Cormode, L. Golab, F. Korn, A. McGregor, D. Srivastava, X. Zhang. Estimating the confidence of conditional functional dependencies, SIGMOD Conference, ACM, 2009, 469–482
- [24] V. Crescenzi, G. Mecca, P. Merialdo. Automatic web information extraction in the roadrunner system, In Revised Papers from the HUMACS, DASWIS, ECOMO, and DAMA on ER 2001 Workshops, 2002, 264–277

-
- [25] V. Crescenzi, D. Qiu, P. Merialdo. A Framework for Learning Web Wrappers from the Crowd, Proceedings of the 22nd WWW Conference, 2013, 261–272
- [26] J. Demetrovics, L. Libkin, I. Muchnik. Functional Dependencies in Relational Databases: a Lattice Point of View, Discrete Applied Mathematics, Volume 40, Issue 2, 1992, 155–185
- [27] M. H. Dunham. Data Mining: Introductory and Advanced Topics, Prentice-Hall, 2002
- [28] W. W. Eckerson. Data Quality and the Bottom Line: Achieving Business Success through a Commitment to High Quality Data, Tech. rep., The Data Warehousing Institute, 2002
- [29] D. W. Embley, W.Y. Mok. Developing XML documents with guaranteed “good” properties, ER 2001, 20th International Conference on Conceptual Modeling, 2001, 426–441
- [30] W. Fan, F. Geerts, X. Jia, A. Kementsietsidis. Conditional functional dependencies for capturing data inconsistencies, ACM Trans. Database Syst., 33(2), 2008
- [31] W. Fan, F. Geerts, L. V. S. Lakshmanan, M. Xiong. Discovering conditional functional dependencies, In ICDE, IEEE, 2009, 1231–1234
- [32] W. Fan, S. Ma, Y. Hu, J. Liu, Y. Wu. Propagating functional dependencies with conditions, PVLDB, 1(1):2008, 391–407
- [33] W. Fan, J. Li, S. Ma, N. Tang, W. Yu. Towards certain fixes with editing rules and master data, PVLDB, 3(1):2010, 173–184
- [34] W. Fan, F. Geerts, J. Li, M. Xiong. Discovering conditional functional dependencies, IEEE Trans. Knowl. Data Eng., 23(5):2011, 683–698
- [35] U. Fayyad, G. Piatetsky-Shapiro, P. Smyth. From Data Mining to Knowledge Discovery in Databases, Advances in knowledge discovery and data mining. American Association for Artificial Intelligence, 1996, 1 – 34
- [36] L. Golab, H. J. Karloff, F. Korn, D. Srivastava, B. Yu. On generating near-optimal tableaux for conditional functional dependencies, PVLDB, 1(1):2008, 376–390
- [37] C. F. Goldfarb. The SGML handbook, Oxford University Press, Inc., 1990
- [38] J. Hereth. Relational Scaling and Databases, Proceedings of the 10th International Conference on Conceptual Structures: Integration and Interfaces LNCS 2393, Springer Verlag, 2002, 62–76

- [39] J. Hereth, G. Stumme, R. Wille, U. Wille. Conceptual Knowledge Discovery - a Human-Centered Approach, *Journal of Applied Artificial Intelligence*, 17:2003, 281–301
- [40] T. W. Hong, K. L. Clark. Using grammatical inference to automate information extraction from the web, 5th European Conference on PKDD, 2001, 216–227
- [41] Y. Huhtala, J. Kärkkäinen, P. Porkka, H. Toivonen. TANE: An efficient algorithm for discovering functional and approximate dependencies, *Comput. J.*, 42:1999, 100–111
- [42] D. Ignatov, K. T. Janosi-Rancz, S. Kuznetsov. Towards a framework for near-duplicate detection in document collections based on closed sets of attributes, *Acta Universitatis Sapientiae, Informatica*, 2, 2009, 215–233
- [43] K. T. Janosi-Rancz, V. Varga, T. Nagy Detecting XML Functional Dependencies through Formal Concept Analysis, 14th ADBIS Conference, LNCS 6295, Springer, 2010, 595–598
- [44] K. T. Janosi-Rancz, V. Varga. A method for mining functional dependencies in relational database design using FCA, *Studia Universitatis Babes-Bolyai, Informatica*, vol. LIII, No. 1, 2008, 17–28
- [45] K.T. Janosi-Rancz, A. Lajos. Semantic Data Extraction, Accepted for The 8th International Conference INTER-ENG 2014, Interdisciplinarity in Engineering, 9 - 10 October 2014, “Petru Maior” University of Tîrgu Mureş, Romania
- [46] K.T. Janosi-Rancz. Finding, Managing and Inforcing CFDs and Ars via a semi-automatic learning strategy, 10th Joint Conference on Mathematics and Computer Science, May 21-25, Cluj Napoca, Romania, *Studia Universitatis Babes-Bolyai, Informatica*, 2014
- [47] K.T. Janosi-Rancz, V. Varga. XML Schema Refinement Through Formal Concept Analysis, *Studia Universitatis Babes-Bolyai, Informatica*, Volume LVII, Number 3, 2012, 49–64
- [48] K. T. Janosi-Rancz, V. Varga, J. Puskas, A Soft Tool for Relational Database Design using Formal Concept Analysis, 7th Joint Conference on Mathematics and Computer Science Cluj, 3-6 July 2008, *Studia Universitatis Babes-Bolyai, Informatica*, vol. LIII, No. 2:2008, 67–78
- [49] M. Kaytoue. Mining numerical data with formal concept analysis and pattern structures (PhD Thesis), Université Henri Poincaré - Nancy, 2011
- [50] A. Kirsch, M. Mitzenmacher, A. Pietracaprina, G. Pucci, E. Upfal, F. Vandin, An efficient rigorous approach for identifying statistically significant frequent itemsets, In

- Proceedings of the 28th ACM SIGMOD-SIGACT-SIGART Symposium on Principles of Database Systems, Providence, Rhode Island, USA, 2009.
- [51] S. Kolahi, L. V. S. Lakshmanan. On approximating optimum repairs for functional dependency violations, *ICDT*, vol. 361 of ACM ICPS, 2009, 53–62
- [52] S. Kolahi, L. Libkin. XML design for relational storage, In *Proceedings of the 16th International World Wide Web Conference*, 2007, 1083–1092
- [53] S.O. Kuznetsov. A fast algorithm for computing all intersections of objects in a finite semilattice, *Automatic documentation and Mathematical linguistics* 27(5), 1993, 11–21
- [54] S.O. Kuznetsov, S. A. Obiedkov. Comparing performance of algorithms for generating concept lattices, *Journal of Experimental and Theoretical Artificial Intelligence*, Vol. 14, Nos 2-3, 2002, 189–216
- [55] A. H. F. Laender, B. A. Ribeiro-Neto, A. S. da Silva, J. S. Teixeira. A brief survey of web data extraction tools, *SIGMOD Rec.* 31, 2/2002, 84–93
- [56] M. Lee, T. Ling, W.L. Low. Designing functional dependencies for XML, *EDBT Conference*, 2002, 124–141
- [57] J. Liu, M. W. Vincent, C. Liu. Local XML functional dependencies, In *Proceedings of the 5th ACM International Workshop on Web Information and Data Management*, 2003, 23–28
- [58] S. Lopes, J. M. Petit, L. Lakhal. Functional and approximate dependency mining: database and fca points of view, *J. Exp. Theor. Artif. Intell.*, 14(2-3):, 2002, 93–114
- [59] R. Medina, L. Nourine. Conditional functional dependencies: An fca point of view, In L. Kwuida and B. Sertkaya, editors, *ICFCA*, volume 5986 of LNCS, Springer, 2010, 161–176
- [60] R. Medina, L. Nourine. A unified hierarchy for functional dependencies, conditional functional dependencies and association rules, In *ICFCA*, 2009, 98–113
- [61] A. Mesbah, A. van Deursen, S. Lenselink. Crawling ajax-based web applications through dynamic analysis of user interface state changes, *ACM Trans. Web*, 6(1), 2012, 1–30
- [62] I. Muslea, S. Minton, C. A. Knoblock. Active learning with multiple views. *J. Artif. Intell. Res. (JAIR)*, 27:2006, 203–233
- [63] S. Perugini, N. Ramakrishnan. Mining web functional dependencies for flexible information access, *J. Am. Soc. Inform. Sci. Technol.*, 58:2007, 1805–1819.

- [64] J. Poelmans, D. I. Ignatov, S. O. Kuznetsov, G. Dedene: Formal concept analysis in knowledge processing: A survey on applications, *Expert Systems with Applications*, 40(16):2013, 6538–6560
- [65] U. Priss. Formal Concept Analysis in Information Science, In Cronin, B. (ed.) *Annual Review of Information Science and Technology*, ASIST, vol. 40, 2008
- [66] K. Probst, R. Ghani, M. Krema, A. E. Fano, Y. Liu. Semi-supervised learning of attribute-value pairs from product descriptions, *IJCAI.*, 2007, 2838–2843
- [67] G. Schreiber, H. Akkermans, A. Anjewierden, R. de Hoog, N. R. Shadbolt, W. van de Velde, B. J. Wielinga. *Knowledge Engineering and Management: The CommonKADS Methodology*, MIT Press, 2000
- [68] I. N. Md. Shaharane, F. Hadzic, T. S. Dillon. Interestingness measures for association rules based on statistical validity, *Knowledge-Based Systems*, Volume 24, Issue 3, 2011, 386–392
- [69] K. Techapichetvanich, A. Datta. VisAR : A new technique for visualizing mined association rules, *ADMA.*, 2005, 88–95
- [70] H. S. Thompson, D. Beech, M. Maloney, N. Mendelsohn (Editors). *XML schema part 1: Structures second edition*, 2004, <http://www.w3.org/tr/xmlschema-1/>
- [71] J. D. Ullman. *Principles of Database Systems*. Computer Science Press, 1982
- [72] P. Valtchev, R. Missaoui, R. Godin, A framework for incremental generation of closed itemsets, *Discrete Applied Mathematics* 156, Elsevier, 2008, 924–949
- [73] V. Varga, K. T. Janosi Rancz. A Software Tool to Transform Relational Databases in Order to Mine Functional Dependencies in it Using Formal Concept Analysis, *Proceedings of the Sixth International Conference on Concept Lattices and Their Applications*, Olomouc, Czech Republic, 2008, 1–9
- [74] V. Varga, K. T. Janosi Rancz, C. Sacarea, K. Csioban, XML Design: an FCA Point of View, *Proceedings of 2010 IEEE International Conference on Automation, Quality and Testing, Robotics*, Theta 17th edition, 2010, 165–170
- [75] M. W. Vincent, J. Liu. Functional dependencies for XML, In *Proceedings of the 5th Asian-Pacific Web Conference*, 2003, 22–34
- [76] M. W. Vincent, J. Liu, C. Liu. Strong functional dependencies and their application to normal forms in XML, *ACM TODS*, 29(3): 2004, 445–462
- [77] M. W. Vincent, J. Liu, M. Mohania, On the Equivalence between FDs in XML and FDs in Relations, *Acta Informatica* 44(3-4), 2007, 207–247

-
- [78] F. Vogt, R. Wille. TOSCANA - a Graphical Tool for Analyzing and Exploring Data, Graph Drawing, 1994, 226-233
- [79] R. Wille. Restructuring lattice theory: an approach based on hierarchies of concepts, I. Rival (Ed.): Ordered sets, Reidel. Dordrecht-Boston, 1982, 445–470
- [80] B. Ganter and R. Wille. Formal Concept Analysis: Mathematical Foundations, Springer, 1999
- [81] R. Wille. Why can concept lattices support knowledge discovery in databases? Journal of Experimental and Theoretical Artificial Intelligence, 14(2-3), 2004, 81–92
- [82] G. I. Webb. Discovering Significant Patterns, Machine Learning, Springer, 2007, 1–33
- [83] C. Wyss, C. Giannella, E. Robertson. FastFDs: A heuristic-driven, depth-first algorithm for mining functional dependencies from relation instances, Proc. of 3rd Int. Conf. of the Data Warehousing and Knowledge Discovery, 2001
- [84] T. L. Wong, W. Lam. An unsupervised method for joint information extraction and feature mining across different Web site, Data And Knowledge Engineering, 68, 2009, 107–125
- [85] T. L. Wong, W. Lam. Unsupervised Extraction of Popular Product Attributes from Web Sites, 8th Asia Information Retrieval Societies Conference, 2012, 437–446
- [86] W3C. XML Schema, 2004, <http://www.w3.org/TR/xmlschema-0/>
- [87] H. Yao, H.J. Hamilton, C.J. Butz, *FD_MINE*: Discovering Functional Dependencies in a Database Using Equivalences, IEEE Int. Conf. on Data Mining (ICDM02), 2002, 9–12
- [88] M. Yakout, A. K. Elmagarmid, J. Neville, M. Ouzzani, and I. F. Ilyas. Guided data repair, PVLDB, 4(5):2011, 279–289
- [89] S. A. Yevtushenko. System of data analysis “Concept Explorer”, Proceedings of the 7th national conference on Artificial Intelligence KII-2000 Russia, 2000, 127–134
- [90] C. Yu, H. V. Jagadish. XML schema refinement through redundancy detection and normalization, VLDB J. 17(2):2008, 203–223
- [91] Y. Zhai, B. Liu. Extracting web data using instance-based learning, World Wide Web, 10, 2007, 113–132
- [92] X. Zheng, Y. Gu, Y. Li. Data extraction from web pages based on structural-semantic entropy, In Proceedings of the 21st WWW Conference, 2012, 93–102